



Europäisches  
Patentamt

European  
Patent Office

Office européen  
des brevets

Bescheinigung

Certificate

Attestation

Die angehefteten Unterlagen stimmen mit der ursprünglich eingereichten Fassung der auf dem nächsten Blatt bezeichneten europäischen Patentanmeldung überein.

The attached documents are exact copies of the European patent application described on the following page, as originally filed.

Les documents fixés à cette attestation sont conformes à la version initialement déposée de la demande de brevet européen spécifiée à la page suivante.

Patentanmeldung Nr. Patent application No. Demande de brevet n°

00202437.0

Der Präsident des Europäischen Patentamts;  
Im Auftrag

For the President of the European Patent Office

Le Président de l'Office européen des brevets  
p.o.

I.L.C. HATTEN-HECKMAN

DEN HAAG, DEN  
THE HAGUE, 24/07/01  
LA HAYE, LE

EPA/EPO/OEB Form 1014 - 02.91

**CERTIFIED COPY OF  
PRIORITY DOCUMENT**

**THIS PAGE BLANK (USPTO)**



Europäisches  
Patentamt

European  
Patent Office

Office européen  
des brevets

**Blatt 2 der Bescheinigung**  
**Sheet 2 of the certificate**  
**Page 2 de l'attestation**

Anmeldung Nr.:  
Application no.:  
Demande n°: 00202437.0

Anmeldetag:  
Date of filing:  
Date de dépôt: 10/07/00 ✓

Anmelder:  
Applicant(s):  
Demandeur(s):  
STMicroelectronics S.r.l.  
20041 Agrate Brianza (Milano)  
ITALY

Bezeichnung der Erfindung:  
Title of the invention:  
Titre de l'invention:  
A method of compressing digital images

In Anspruch genommene Priorität(en) / Priority(ies) claimed / Priorité(s) revendiquée(s)

Staat:  
State:  
Pays:

Tag:  
Date:  
Date:

Aktenzeichen:  
File no.  
Numéro de dépôt:

Internationale Patentklassifikation:  
International Patent classification:  
Classification internationale des brevets:

H04N7/30

Am Anmeldetag benannte Vertragsstaaten:  
Contracting states designated at date of filing: AT/BE/CH/CY/DE/DK/ES/FI/FR/GB/GR/IE/IT/LI/LU/MC/NL/PT/SE/TR  
Etats contractants désignés lors du dépôt:

Bemerkungen:  
Remarks:  
Remarques:

**THIS PAGE BLANK (USPTO)**

10. 07. 2000

(41)

1

The present invention relates to a method of compressing digital images.

Digital images are commonly used in several applications such as, for example, in digital still cameras (DSC). A digital image consists of a matrix of elements, commonly referred to as a bit map; each element of the matrix, which represents an elemental area of the image (a pixel or pel), is formed by digital values indicating corresponding components of the pixel.

Digital images are typically subjected to a compression process in order to increase the number of digital images which can be stored simultaneously, such as onto a memory of the camera; moreover, this allows transmission of digital images (for example in the INTERNET) to be easier and less time consuming. A compression method commonly used in standard applications is the JPEG (Joint Photographic Experts Group) algorithm, described in CCITT T.81, 1992.

In the JPEG algorithm, 8x8 pixel blocks are extracted from the digital image; Discrete Cosine Transform (DCT) coefficients are then calculated for the components of each block. The DCT coefficients are rounded off using corresponding quantization tables; the quantized DCT coefficients are encoded in order to obtain a compressed digital image (from which the corresponding

original digital image can be extracted later on by a decompression process).

In some applications, it is necessary to provide a substantially constant memory requirement for each compressed digital image (the so called Compression Factor Control, or CF-CTRL). This problem is particularly perceived in digital still cameras; in fact, in this case it must be ensured that a minimum number of compressed digital images can be stored onto the memory of the camera, in order to guarantee that a minimum number of photos can be taken by the camera.

The compression factor control is quite difficult in algorithms, such as the JPEG, wherein the size of the compressed digital image depends on the content of the corresponding original digital image.

Generally, the compression factor is controlled by scaling the quantization tables using a multiplier coefficient (gain factor). The gain factor to obtain a target compression factor is determined using iterative methods. The compression process is executed several times, at least twice; the gain factor is modified according to the result of the preceding compression process, until the compressed digital image has a size that meets the target compression factor.

The methods known in the art require a high

computation time, so that they are quite slow. Moreover, the known methods involve a considerable power consumption; this drawback is particular acute when the compression method is implemented in a digital still  
5 camera, or other portable devices which are supplied by batteries.

It is an object of the present invention to overcome the above mentioned drawbacks. In order to achieve this object, a method of compressing a digital image as set  
10 out in the first claim is proposed.

Briefly, the present invention provides a method of compressing a digital image including a matrix of elements each one consisting of at least one component of different type representing a pixel, the method  
15 comprising the steps of splitting the digital image into a plurality of blocks and calculating, for each block, a group of DCT coefficients for the components of each type, and quantizing the DCT coefficients of each group using a corresponding quantization table scaled by a gain  
20 factor for achieving a target compression factor; the method also comprises the steps of further quantizing the DCT coefficients of each group using the corresponding quantization table scaled by a pre-set factor, arranging the further quantized DCT coefficients in a zig-zig  
25 vector, calculating a basic compression factor provided

by the quantization table scaled by the pre-set factor as a first function of the zigzag vector, and estimating the gain factor as a second function of the basic compression factor, the second function being determined  
5 experimentally according to the target compression factor.

Moreover, the present invention also provides a corresponding device for compressing a digital image and a digital still camera comprising this device.

10 Further features and the advantages of the solution according to the present invention will be made clear by the following description of a preferred embodiment thereof, given purely by way of a non-restrictive indication, with reference to the attached figures, in  
15 which:

Fig.1 is a schematic block diagram of a digital still camera, in which the compression method of the invention can be used,

Figg.2 depicts an example of relation basic  
20 compression factor/gain factor,

Figg.3a-3b show a flow chart of the compression method.

With reference in particular to Fig.1, this shows a digital still camera 100 for taking digital images  
25 representative of real scenes. A digital image is



5

constituted by a matrix with N rows and M columns (for example, 640 rows by 480 columns); each element of the matrix consists of one or more digital values (for example three values each one of 8 bits, ranging from 0 to 255) representative of respective optical components of a pixel.

The camera 100 includes an image-acquisition unit 105 formed by a diaphragm and a set of lenses for transmitting the light corresponding to the image of the real scene onto a sensor unit (SENS) 110. The sensor unit 110 is typically constituted by a Charge-Coupled Device (CCD); a CCD is an integrated circuit which contains a matrix of light-sensitive cells, each one generating a voltage the intensity of which is proportional to the exposure of the light-sensitive cell. The voltage generated by each light-sensitive cell is supplied to an analog/digital converter, which produces a corresponding digital value.

In order to reduce the number of light-sensitive cells, the sensor unit 110 does not detect all the components for every pixel; typically, only one light-sensitive cell is provided for each pixel. The CCD is covered by a colour filter consisting of a matrix of filter elements each one associated with a corresponding light-sensitive cell of the CCD; each filter element transmits (absorbing a minimal portion) the luminous

6

radiation belonging only to the wavelength of red, blue or green light (substantially absorbing the others), so as to detect a red color component (R), a green color component (G), or a blue color component (B) for each  
5 pixel.

In particular, the filter is of the Bayer type as described in US-A-3,971,065, in which only the G component is detected for a half of the pixels, in a chessboard-like arrangement; the R component or the B  
10 component is detected for the other half of the pixels, in respective alternate rows, as shown in the following table:

	...	...	...	...	...	...	...	...	...	...
	...	G	R	G	R	G	R	G	R	G
15	...	B	G	B	G	B	G	B	G	B
	...	G	R	G	R	G	R	G	R	G
	...	B	G	B	G	B	G	B	G	B
	...	...	...	...	...	...	...	...	...	...

An incomplete digital image SImg, in which each element  
20 consists of a single colour component (R, G or B), is output by the sensor unit 110.

The camera 100 includes a control unit 115 formed by several blocks which are connected in parallel to a communication bus 120. Particularly, a pre-processing  
25 unit (PRE\_PROC) 125 receives the incomplete digital image

SImg. The pre-processing unit 125 determines various parameters of the incomplete digital image SImg (such as a high-frequency content and an average luminosity); these parameters are used to automatically control a focus (auto-focus) and an exposure (auto-exposure) by means of corresponding control signals Sc which are supplied to the acquisition unit 105. The pre-processing unit 125 also modifies the incomplete digital image SImg, for example applying a white-balance algorithm which corrects the colour shift of the light towards red (reddish) or towards blue (bluish), in dependence on the colour temperature of the light source; a corresponding incomplete digital image BImg is output by the pre-processing unit 125 and sent onto the bus 120.

The incomplete digital image BImg is received by an image-processing unit (IPU) 130. The image-processing unit 130 interpolates the missing colour components in each element of the incomplete digital image BImg, in order to obtain a corresponding digital image RGB wherein each pixel is represented by the R component, the G component and the B component. The digital image RGB is then processed to improve image quality, for example correcting exposure problems such as back-lighting or excessive front illumination, reducing a noise introduced by the CDD, correcting alterations of a selected colour

tone, applying special effects (such as a mist effect), compensating the loss of sharpness due to a  $\gamma$ -correction function (typically applied by a television set); moreover, the digital image can be enlarged, a particular  
5 of the image can be zoomed, or the ratio of its dimensions can be changed (for example from 4:3 to 16:9), and the like.

The digital image RGB is then converted into a corresponding digital image YUV in a  
10 luminance/chrominance space. Each pixel of the digital image YUV is represented by a luminance component Y (providing information about the brightness), and two chrominance components Cu and Cv (providing information about the hue); the Y,Cu,Cv components are calculated  
15 from the respective R,G,B components applying, for example, the following equations:

$$Y = 0.299 \cdot R + 0.587 \cdot G + 0.114 \cdot B$$

$$Cu = -0.1687 \cdot R - 0.3313 \cdot G + 0.5 \cdot B + 128$$

$$Cv = 0.5 \cdot R - 0.4187 \cdot G - 0.0813 \cdot B + 128$$

20 This allows chrominance information to be easily identified, in order to discard more chrominance information than luminance information during a following compression process of the digital image (being the human eye more sensitive to luminance than chrominance). The  
25 digital image YUV is sent onto the bus 120.

A compression unit 135 is also connected to the bus 120; the compression unit 135 receives the digital image YUV and outputs a corresponding digital image JImg compressed applying a JPEG algorithm. The compression  
 5 unit 135 includes a Discrete Cosine Transform (DCT) unit 145, which is input the digital image YUV. Each component of the digital image YUV is shifted from the range 0..255 to the range -128..+127, in order to normalize the result of the operation. The digital image YUV is then split  
 10 into several blocks of 8x8 pixels (640x480/64 = 4800 blocks in the example at issue). Each block of Y components BLy, each block of Cu components BLu, and each block of Cv components BLv is translated into a group of DCT coefficients DCTy, a group of DCT coefficients DCTu,  
 15 and a group of DCT coefficients DCTv, respectively, representing a spatial frequency of the corresponding components. The DCT coefficients  $DCT_{y,u,v}[h,k]$  (with  $h,k=0..7$ ) are calculated using the following formula:

$$DCT_{y,u,v}[h,k] = \frac{1}{4} Dh Dk \sum_{x=0}^7 \sum_{y=0}^7 BL_{y,u,v}[x,y] \cos \frac{(2h+1)x\pi}{16} \cos \frac{(2k+1)y\pi}{16}$$

wherein  $Dh, Dk = 1/\sqrt{2}$  for  $h,k=0$  and  $Dh, Dk=1$  otherwise. The  
 20 first DCT coefficient of each group is referred to as DC coefficient, and it is proportional to the average of the components of the group, whereas the other DCT coefficients are referred to as AC coefficients. The

10

groups of DCT coefficients  $DCT_{y,u,v}$  are sent onto the bus 120.

The compression unit 135 further includes a quantizer (QUANT) 150, which receives (from the bus 120) the groups of DCT coefficients  $DCT_{y,u,v}$  and a scaled quantization table for each type of component; typically, a scaled quantization table  $SQ_y$  is used for the Y components and a scaled quantization table  $SQ_{uv}$  is used for both the  $C_u$  components and the  $C_v$  components. Each scaled quantization table  $SQ_y, SQ_{uv}$  consists of a  $8 \times 8$  matrix of quantization constants; the DCT coefficients of each group are divided by the corresponding quantization constants and rounded off to the nearest integer. As a consequence, smaller and unimportant DCT coefficients disappear and larger DCT coefficients lose unnecessary precision. The quantization process generates corresponding groups of quantized DCT coefficients  $QDCT_y$  for the Y component, groups of quantized DCT coefficients  $QDCT_u$  for the  $C_u$  component, and groups of quantized DCT coefficients  $QDCT_v$  for the  $C_v$  component.

These values drastically reduce the amount of information required to represent the digital image. The JPEG algorithm is then a lossy compression method, wherein some information about the original image is finally lost during the compression process; however, no

11

image degradation is usually visible to the human eye at normal magnification in the corresponding de-compressed digital image for a compression ratio ranging from 10:1 to 20:1 (defined as the ratio between the number of bits  
5 required to represent the digital image YUV and the number of bits required to represent the compressed digital image JImg).

Each scaled quantization table  $SQ_y, SQ_{uv}$  is obtained multiplying a corresponding quantization table  $Q_y, Q_{uv}$  by  
10 a gain factor  $G$  (determined as set out in the following), that is  $SQ_y = G \cdot Q_y$  and  $SQ_{uv} = G \cdot Q_{uv}$ . The gain factor  $G$  is used to obtain a desired, target compression factor  $bp_t$  of the JPEG algorithm (defined as the ratio between the number  
15 of bits of the compressed digital image JImg and the number of pixels). Particularly, if the gain factor  $G$  is greater than 1, the compression factor is reduced (compared to the one provided by the quantization tables  
20  $Q_y, Q_{uv}$ ), whereas if the gain factor  $G$  is less than 1 the compression factor is increased.

20 The quantization tables  $Q_y, Q_{uv}$  are defined so as to discard more chrominance information than luminance information. For example, the quantization table  $Q_y$  is:

	1	11	10	16	24	40	51	61
	12	12	14	19	26	58	60	55
25	14	13	16	24	40	57	69	56

12

	14	17	22	29	51	87	8	62
	18	22	37	56	68	109	203	77
	24	35	55	64	81	104	113	92
	49	64	78	87	103	121	120	101
5	72	92	95	98	112	100	103	99

and the quantization table Quv is:

	1	18	24	47	99	99	99	99
	18	21	26	66	99	99	99	99
	24	26	56	99	99	99	99	99
10	47	66	99	99	99	99	99	99
	99	66	99	99	99	99	99	99
	99	66	99	99	99	99	99	99
	99	66	99	99	99	99	99	99
	99	66	99	99	99	99	99	99

15 Preferably, the quantization constants for the DC coefficients are equal to 1 in both cases, in order not to loose any information about the mean content of each block, and then to avoid the so-called "block-effect" (wherein a contrast is perceivable between the blocks of

20 the de-compressed image).

The groups of quantized DCT coefficients QDCTy,u,v are directly provided to a zigzag unit (ZZ) 155. The zigzag unit 155 modifies and reorders the quantized DCT coefficients to obtain a single vector ZZ of digital

25 values. Each quantized DC coefficient (but the one of a



first group) is represented as the difference from the quantized DC coefficient of a previous group. The quantized AC coefficients are arranged in a zigzag order, so that quantized AC coefficients representing low frequencies are moved to the beginning of the group and quantized AC coefficients representing high frequencies are moved to the end of the group; since the quantized AC coefficients representing high frequencies are more likely to be zeros, this increases the probability of having longer sequences of zeros in the vector ZZ (which require a lower number of bits in a run length encoding scheme).

The vector ZZ is directly provided to an encoder (ENC) 160, which also receives one or more encoding tables HT from the bus 120. Each value of the vector ZZ is encoded using a Huffman scheme, wherein the value is represented by a variable number of bits which is inversely proportional to a statistical frequency of use thereof. The encoder 160 then generates the corresponding compressed digital image JImg (which is sent onto the bus 120). The compressed digital image JImg is typically formed by a header (for example some tens of bytes containing information about the digital image and the compression method, such as the quantization tables and the dimension of the digital image) followed by the

encoded values. If the last encoded value associated with a block is equal to 00, it must be followed by a (variable) End of Block (EOB) control word. Moreover, if an encoded value is equal to a further control word FF  
5 (used as a marker), this value must be followed by a 00 value.

The control unit 115 also includes a working memory 165, typically a SDRAM (Synchronous Dynamic Random Access Memory) and a microprocessor ( $\mu$ P) 170, which controls the  
10 operation of the device. Several peripheral units are further connected to the bus 120 (by means of a respective interface). Particularly, a non-volatile memory 175, typically a flash E<sup>2</sup>PROM, stores the quantization tables Qy, Quv, the encoding tables HT, and a  
15 control program for the microprocessor 170. A memory card (MEM\_CARD) 180 is used to store the compressed digital images JImg; the memory card 185 has a capacity of a few Mbytes, and can store several tens of compressed digital images JImg. At the end, the camera 100 includes an  
20 input/output (I/O) unit 185 consisting, for example, of a series of push-buttons, for enabling the user to select various functions of the camera 100 (such as an on/off button, an image quality selection button, a shot button, a zoom control button), and a liquid-crystal display  
25 (LCD), for supplying data on the operative state of the camera 100 to the user.

Likewise considerations apply if the camera has a different architecture or includes different units, such as equivalent communication means, a CMOS sensor, a viewfinder or an interface for connection to a personal computer (PC) and a television set, if another colour filter (not with a Bayer pattern) is used, if the compressed digital images are directly sent outside the camera (without being stored onto the memory card), and so on. Alternatively, the digital image includes a single gray component for each element, the digital image is converted into another space (not a luminance/chrominance space), the digital image RGB is directly compressed (without being converted), the digital image YUV is manipulated to down-sample the  $C_u, C_v$  components by averaging groups of pixels together (in order to eliminate further information without sacrificing overall image quality), or no elaboration of the digital image is performed; similarly, one or more different quantization tables are used, arithmetic encoding schemes are employed, a different compression algorithm is used (such as a progressive JPEG). Moreover, the compression method of the present invention leads itself to be implemented even in a different apparatus, such as a portable scanner, a computer in which graphic applications are provided, and the like.

The inventors have discovered that the gain factor  $G$  for obtaining the target compression factor  $bp_t$  is a function of a basic compression factor  $bp_b$ , obtained using the quantization tables  $Q_y, Q_{uv}$  scaled by a pre-set factor  $S$  (determined as set out in the following). The function depends on the target compression factor  $bp_t$  (in addition to the characteristics of the camera 100, such as the dimension of the CCD, the size of the digital image, the quantization tables used), and can be determined a priori by a statistical analysis.

For example, Fig.2 shows a relation between the basic compression factor  $bp_b$  and the gain factor  $G$  for a camera having a CDD with 1 million of light-sensitive cells and for images of 640x480 pixels, with a factor  $S=0,2$  and a target compression factor  $bp_t=2$  bit/pel. This relation can be interpolated as a quadratic function; in other words, the gain factor  $G$  can be estimated using the relation  $G=C_2 \cdot bp_b^2 + C_1 \cdot bp_b + C_0$  (wherein  $C_2$ ,  $C_1$  and  $C_0$  are parameters depending on the characteristics of the camera 100 and the target compression factor  $bp_t$ ).

In order to calculate the basic compression factor  $bp_b$ , the quantizer 150 is supplied with scaled quantization tables  $SQ_y, SQ_{uv}$  obtained multiplying the corresponding quantization tables  $Q_y, Q_{uv}$  by the pre-set factor  $S$ , that is  $SQ_y=S \cdot Q_y$  and  $SQ_{uv}=S \cdot Q_{uv}$ . The quantizer

17

150 determines the corresponding groups of quantized DCT coefficients QDCT<sub>y,u,v</sub> and the zigzag unit 155 modifies and reorders the quantized DCT coefficients to obtain the vector ZZ.

5       The vector ZZ is directly provided to a counting unit (COUNT) 190, which outputs the number of bits ZZbits required to encode (in the compressed digital image JImg) the values of the vector ZZ associated with each block. To this end, a look-up table JN is stored onto the E<sup>2</sup>PROM  
10 175 and it is sent to the counting unit 190 (by means of the bus 120); each row of the look-up table JN, addressable by the values of the vector ZZ associated with a block, contains the respective number ZZbits.

      The basic compression factor  $bp_b$  is calculated  
15 summing the numbers ZZbits associated with every block. A constant value indicating the number of bits required to encode the header of the compressed digital image JImg is then added to the sum. The result is divided by the number of pixels (N·M).

20       More generally, the method of the present invention includes the steps of further quantizing the DCT coefficients of each group using the corresponding quantization table scaled by a pre-set factor, arranging the further quantized DCT coefficients in a zig-zig  
25 vector, calculating a basic compression factor provided

by the quantization table scaled by the pre-set factor as a first function of the zigzag vector, and estimating the gain factor as a second function of the basic compression factor, the second function being determined  
5 experimentally according to the target compression factor.

The method of the invention is very fast, in that only some of the operations performed by the compression unit (i.e., the quantization and the zigzag reordering)  
10 are executed twice. In this respect, it should be noted that the operations performed by the counting unit 190 are far simpler and faster than the ones performed by the encoder 160.

The solution according to the present invention is  
15 particularly advantageous in portable devices supplied by batteries (even if different applications are not excluded), since it drastically reduces the power consumption.

These results are achieved with a low error (of the  
20 order of a few units per cent) between the target compression factor  $bp_t$  and a compression factor  $bp_a$  actually obtained, defined as  $(bp_t - bp_a)/bp_t$ . Experimental results on the camera at issue provided a mean error of -1% (the negative error is more important than the  
25 positive error because the size of the compressed digital

image is bigger than the target one), with a distribution of 98% between  $\pm 6\%$  and 100% between  $\pm 10\%$ .

In the above described architecture, a single quantizer 150 is provided. The quantizer 150 is supplied with the scaled quantization tables  $SQ_y, SQ_{uv}$  obtained multiplying the corresponding quantization tables  $Q_y, Q_{uv}$  by the pre-set factor  $S$  (for calculating the number  $ZZ_{bits}$ ) or with the scaled quantization tables  $SQ_y, SQ_{uv}$  obtained multiplying the corresponding quantization tables  $Q_y, Q_{uv}$  by the gain factor  $G$  (for generating the compressed digital image  $JImg$ ). This solution is particularly simple and flexible.

Preferably, two or more sets of parameters  $C_2, C_1, C_0$ , each one associated with a different value of the target compression factor  $bp_t$  and with a different size of the digital image, are determined a priori by a statistical analysis. A look-up table, wherein each row addressable by the value of the target compression factor  $bp_t$  contains the respective parameters  $C_2, C_1, C_0$ , is stored onto the  $E^2PROM$  175. This feature allows different compression factors to be easily selected by the user.

Advantageously, the factor  $S$  is determined a priori by a statistical analysis, in order to further reduce the error between the target compression factor  $bp_t$  and the actual compression factor  $bp_a$ . Experimental results have

shown that the factor  $S$  which minimizes the error also depends on the target compression factor  $bp_t$  (in addition to the characteristics of the camera 100).

Alternatively, the basic compression factor  $bp_b$  is  
5 calculated in a different manner (for example by software from the whole vector  $ZZ$ ), the relation  $bp_b/E$  is interpolated with a different function (such as a logarithmic function), the look-up table with the parameters  $C2, C1, C0$  is stored elsewhere or a different  
10 memory structure is used, the tables  $Qy, Quv$  are embedded in the quantizer 150 (which is supplied with the pre-set factor  $S$  or the gain factor  $G$ ), or more generally the quantizer is operated in two different conditions (using the quantization tables scaled by the pre-set factor  $S$  or  
15 the gain factor  $G$ , respectively), and the like. However, the method of the present invention leads itself to be carried out even with two distinct quantizers, with only one set of parameters  $C2, C1, C0$ , with the quadratic function implemented by software, with the factor  $S$  set  
20 to a constant value, even equal to 1 (irrespective of the target compression factor  $bp_t$ ).

In order to explain the operation of the camera, reference is made to Figg.3a-3b (together with Fig.1). When the camera 100 is switched on by the user (acting on  
25 the on/off button), the microprocessor 170 runs the



21

control program stored in the E<sup>2</sup>PROM 175. A method 300 corresponding to this control program starts at block 305 and then passes to block 310, wherein the user selects the desired quality of the image (such as low or high) by  
5 acting on the corresponding button; the microprocessor 170 determines and stores onto the SDARM 165 the target compression factor  $bp_t$  corresponding to the selected image quality (for example, 1 bit/pel for the low quality and 2 bit/pel for the high quality).

10 The method checks at block 315 if the shot button has been partially pressed in order to focus the image; if not, the method returns to block 310; as soon as the user partially presses the shot button, the method proceeds to block 320, wherein the incomplete digital  
15 image SImg is acquired by the sensor unit 110 (the diaphragm is always open and the light is focused by the lenses, through the Bayer filter, onto the CCD). The pre-processing unit 125 then controls the acquisition unit 115 (by means of the control signals  $Sc$ ) according to the  
20 content of the incomplete digital image SImg.

The method checks again the status of the shot button at block 325. If the shot button has been released, the method returns to block 310, whereas if the shot button has been completely pressed (in order to take  
25 a photo) the method continues to block 330; on the other

hand, if no action is performed by the user, the method stays in block 325 in an idle loop.

Considering now block 330, the incomplete digital image SImg is acquired by the sensor unit 110 and  
5 modified by the pre-processing unit 125; the corresponding incomplete digital image BImg is stored onto the SDRAM 165. The method passes to block 335, wherein the incomplete digital image BImg is read from the SDRAM 165 and provided to the image-processing unit  
10 130. The image-processing unit 130 interpolates the missing colour components in each element of the incomplete digital image BImg, in order to obtain the corresponding digital image RGB, and modifies the digital image RGB to improve the image quality; the digital image  
15 RGB is then converted into the corresponding digital image YUV. Proceeding to block 340, the digital image YUV is provided to the DCT unit 140; the DCT unit 140 calculates the groups of DCT coefficients DCT<sub>y,u,v</sub>, which are sent onto the bus 120.

20 The method then forks into two branches which are executed concurrently. A first branch consists of block 345, and a second branch consists of blocks 350-375; the two branches joint at block 378 (described in the following).

25 Considering now block 345, the groups of DCT

coefficients  $DCT_{y,u,v}$  are received and stored onto the SDRAM 165. At the same time, at block 350, the groups of DCT coefficients  $DCT_{y,u,v}$  are also received by the quantizer 150; in the meanwhile, the microprocessor 170  
5 reads the quantization tables  $Q_y, Q_{uv}$  from the  $E^2$ PROM 175 and calculates the scaled quantization tables  $SQ_y, SQ_{uv}$  multiplying the respective quantization tables  $Q_y, Q_{uv}$  by the pre-set factor  $S$ ; the scaled quantization tables  $SQ_y, SQ_{uv}$  are then provided to the quantizer 150.  
10 Continuing to block 355, the quantizer 150 generates the corresponding groups of quantized DCT coefficients  $QDCT_{y,u,v}$ . The method proceeds to block 360, wherein the quantized DCT coefficients  $QDCT_{y,u,v}$  are transformed into the vector  $ZZ$  by the zigzag unit 155.

15 Considering now block 365, the vector  $ZZ$  is provided to the counting unit 190; at the same time, the look-up table  $JN$  is read from the  $E^2$ PROM 175 and sent to the counting unit 190, which determines the number  $ZZ_{bits}$ . The microprocessor 170 receives the number  $ZZ_{bits}$  at  
20 block 370, and calculates the basic compression factor  $bp_b$  accordingly. Continuing now to block 375, the microprocessor reads the parameters  $C_2, C_1, C_0$  associated with the target compression factor  $bp_t$  from the  $E^2$ PROM 175 (addressing the look-up table by the value of the  
25 target compression factor  $bp_t$ ); the microprocessor 170

24

then estimates the gain factor  $G$  for obtaining the target compression factor  $bp_t$  using the read parameters  $C2, C1, C0$ .

Considering now block 378, the groups of DCT  
5 coefficients  $DCTy, u, v$  are read from the SDRAM 165 and sent onto the bus 120. The groups of DCT coefficients  $DCTy, u, v$  are received by the quantizer 150 at block 350a; in the meanwhile, the microprocessor 170 reads the quantization tables  $Qy, Quv$  from the  $E^2$ PROM 175 and  
10 calculates the scaled quantization tables  $SQy, SQuv$  multiplying the respective quantization tables  $Qy, Quv$  by the gain factor  $G$ ; the scaled quantization tables  $SQy, SQuv$  are then provided to the quantizer 150. Continuing to block 355a, the quantizer 150 generates the  
15 corresponding groups of quantized DCT coefficients  $QDCTy, u, v$ . The method proceeds to block 360a, wherein the quantized DCT coefficients  $QDCTy, u, v$  are transformed into the vector  $ZZ$  by the zigzag unit 155. The vector  $ZZ$  is supplied, at block 380, to the encoder 160, which  
20 generates the corresponding compressed digital image  $JImg$ ; the compressed digital image  $JImg$  is then stored onto the SDRAM 165. Continuing to block 385, the compressed digital image  $JImg$  is read from the SDRAM 165 and sent to the memory card 180.

25 The method then checks at block 390 if a stop

25

condition has occurred, for example if the user has switched off the camera 100 (acting on the on/off button) or if the memory card 180 is full. If not, the method returns to block 310; on the other end, the method ends  
5 at block 395.

The preferred embodiment of the present invention described above, with the counting unit implemented in hardware and the basic compression factor calculation and gain factor estimation functions implemented in software,  
10 is a good trade-off between speed and flexibility.

Moreover, this solution requires the operations performed by the DCT unit 145 to be carried out only once.

Likewise considerations apply if the program  
15 executes a different equivalent method, for example with error routines, with sequential processes, and the like. In any case, the method of the present invention leads itself to be carried out even with all the functions completely implemented in hardware or in software, and  
20 with the DCT coefficients calculated twice.

Naturally, in order to satisfy local and specific requirements, a person skilled in the art may apply to the solution described above many modifications and alterations all of which, however, are included within

26

the scope of protection of the invention as defined by  
the following claims.

CLAIMS

(41)

1. A method (300) of compressing a digital image including a matrix of elements each one consisting of at least one component of different type representing a pixel, the method comprising the steps of:

splitting (340) the digital image into a plurality of blocks and calculating, for each block, a group of DCT coefficients for the components of each type,

10 quantizing (350a-355a) the DCT coefficients of each group using a corresponding quantization table scaled by a gain factor for achieving a target compression factor,

characterized by the steps of

further quantizing (350-355) the DCT coefficients of each group using the corresponding quantization table scaled by a pre-set factor,

arranging (360) the further quantized DCT coefficients in a zig-zig vector,

calculating (365-370) a basic compression factor provided by the quantization table scaled by the pre-set factor as a first function of the zigzag vector,

estimating (375) the gain factor as a second function of the basic compression factor, the second function being determined experimentally according to the target compression factor.

25

2. The method (300) according to claim 1, wherein the step (365-370) of calculating the basic compression factor includes the steps of:

determining (365) a first number of bits required to  
5 encode the zigzag vector,

calculating (370) the basic compression factor summing the first number of bits with a second number of bits required to encode control values, and diving the sum by the number of elements of the digital image.

10 3. The method (300) according to claim 1 or 2, wherein the second function is a quadratic function.

4. The method (300) according to any claim from 1 to 3, further comprising the steps of:

storing a plurality of sets of parameters  
15 representing the second function, each set of parameters being associated with a corresponding value of the target compression factor,

selecting (310) an image quality and determining a current value of the target compression factor as a  
20 function of the selected image quality,

reading (375) the parameters associated with the current value of the target compression factor and estimating the gain factor using the read parameters.

5. The method (300) according to any claim from 1 to  
25 4, wherein the pre-set factor is determined



experimentally according to the target compression factor.

6. The method (400) according to any claim from 1 to 5, further comprising the steps of:

5 storing (345) the DCT coefficients onto a working memory and concurrently performing the steps of quantizing (350-355) the DCT coefficients of each group using the corresponding quantization table scaled by the pre-set factor, arranging (360) the quantized DCT  
10 coefficients in the zig-zig vector, calculating (365-370) the basic compression factor, and estimating (375) the gain factor,

reading (378) the DCT coefficients from the working memory for performing the step of quantizing (350a-355a)  
15 the DCT coefficients of each group using the corresponding quantization table scaled by the gain factor.

7. A device (115) for compressing a digital image including a matrix of elements each one consisting of at  
20 least one digital component of different type representing a pixel, the device (115) comprising means (145) for splitting the digital image into a plurality of blocks and calculating, for each block, a group of DCT coefficients for the components of each type, means (150)  
25 for quantizing the DCT coefficients of each group using a

30

corresponding quantization table scaled by a gain factor  
for achieving a target compression factor,

characterized in that

the device (115) further includes means (150) for further  
5 quantizing the DCT coefficients of each group using the  
corresponding quantization table scaled by a pre-set  
factor, means (155) for arranging the further quantized  
DCT coefficients in a zig-zig vector, means (170,190) for  
calculating a basic compression factor provided by the  
10 quantization table scaled by the pre-set factor as a  
first function of the zigzag vector, and means (170) for  
estimating the gain factor as a second function of the  
basic compression factor, the second function being  
determined experimentally according to the target  
15 compression factor.

8. The device (115) according to claim 7, further  
comprising a quantization unit (150) which quantizes the  
DCT coefficients of each group using the corresponding  
quantization table scaled by the gain factor in a first  
20 operative condition and which quantizes the DCT  
coefficients of each group using the corresponding  
quantization table scaled by the pre-set factor in a  
second operative condition.

9. The device (115) according to claim 7 or 8,  
25 wherein the means (170,190) for calculating the basic

compression factor includes means (190) for determining a first number of bits required to encode the zigzag vector, and means (170) for calculating the basic compression factor summing the first number of bits with  
5 a second number of bits required to encode control values and dividing the sum by the number of elements of the digital image

10. The device (115) according to claim 9, further comprising a DCT unit (145) comprising the means for  
10 splitting the digital image and for calculating the DCT coefficients, a zigzag unit (155) comprising the means for arranging the further quantized DCT coefficients in the zig-zig vector, a memory unit (175) for storing the quantization tables, a counting unit (190) comprising the  
15 means for calculating the first number of bits, a processor unit (170) for controlling the device (115), communication means (120) for connecting the DCT unit, the quantization unit, the zigzag unit, the memory unit, the counting unit and the processor unit therebetween,  
20 the processor unit (170) calculating the basic compression factor and estimating the gain factor under the control of a program stored onto the memory unit (175).

11. A digital still camera (100) comprising the  
25 device (115) of any claim from 7 to 10.

**THIS PAGE BLANK (USPTO)**

EPO - DG 1

10. 07. 2000

32

ABSTRACT

(41)

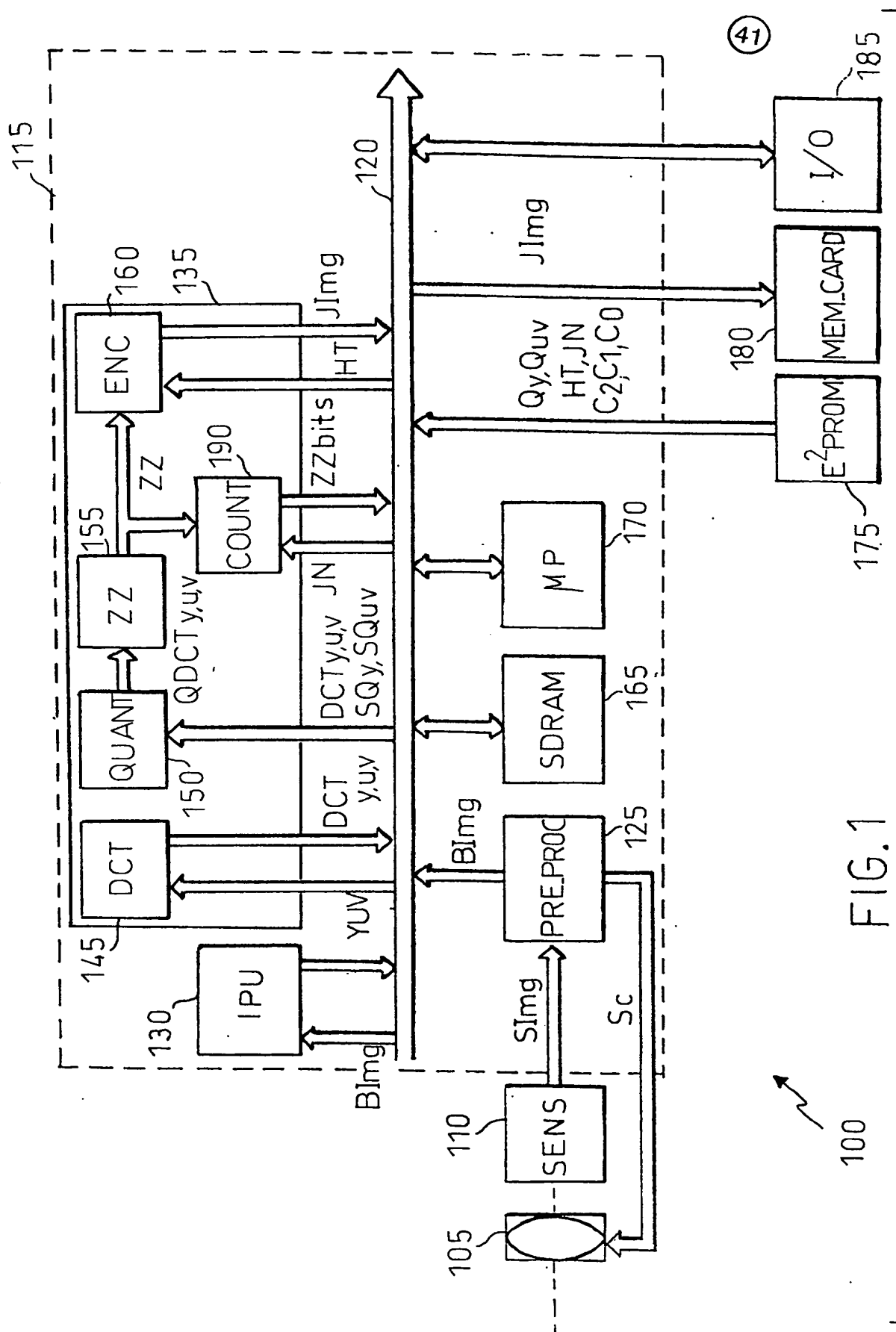
## A METHOD OF COMPRESSING DIGITAL IMAGES

5           A method (300) of compressing a digital image including a matrix of elements each one consisting of at least one component of different type representing a pixel, the method comprising the steps of splitting (340) the digital image into a plurality of blocks and  
10   calculating, for each block, a group of DCT coefficients for the components of each type, and quantizing (350a-355a) the DCT coefficients of each group using a corresponding quantization table scaled by a gain factor for achieving a target compression factor; the method  
15   also comprises the steps of further quantizing (350-355) the DCT coefficients of each group using the corresponding quantization table scaled by a pre-set factor, arranging (360) the further quantized DCT coefficients in a zig-zig vector, calculating (365-370) a  
20   basic compression factor provided by the quantization table scaled by the pre-set factor as a first function of the zigzag vector, and estimating (375) the gain factor as a second function of the basic compression factor, the second function being determined experimentally according  
25   to the target compression factor.

**THIS PAGE BLANK (USPTO)**

10. 07. 2000

1/4



2/4

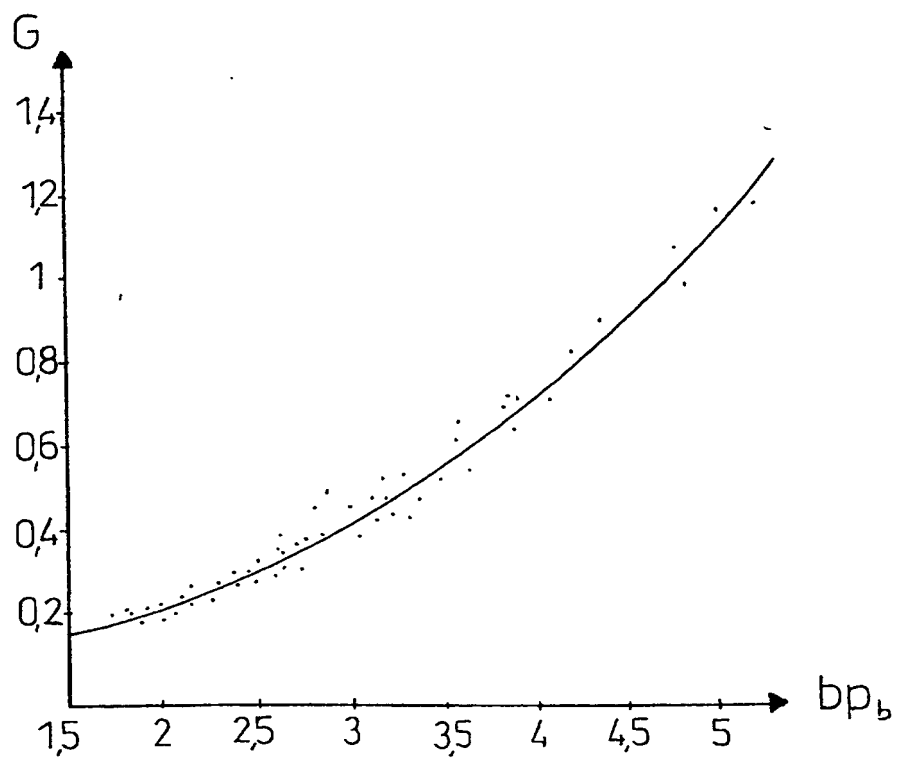
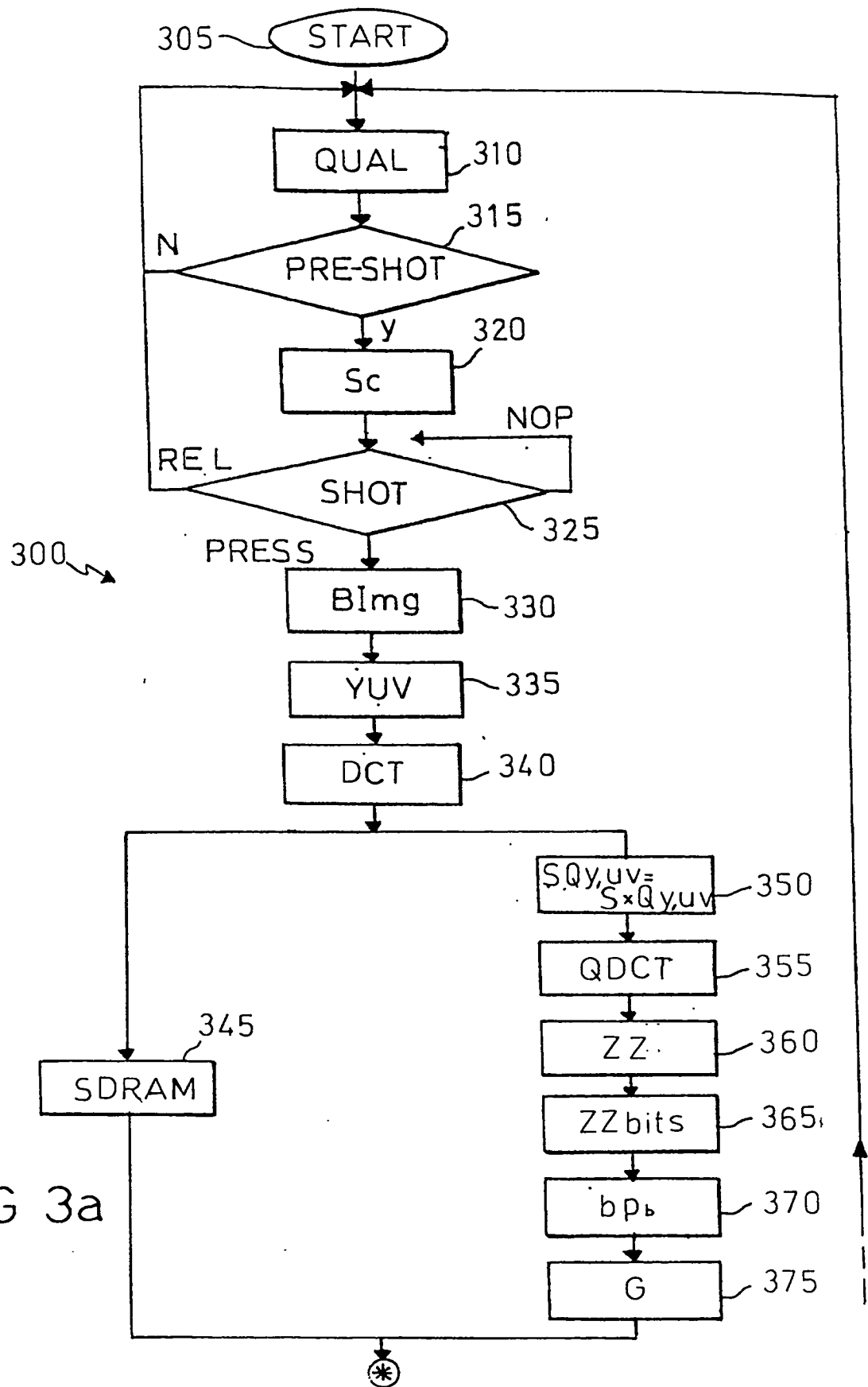


FIG. 2



3/4



4/4

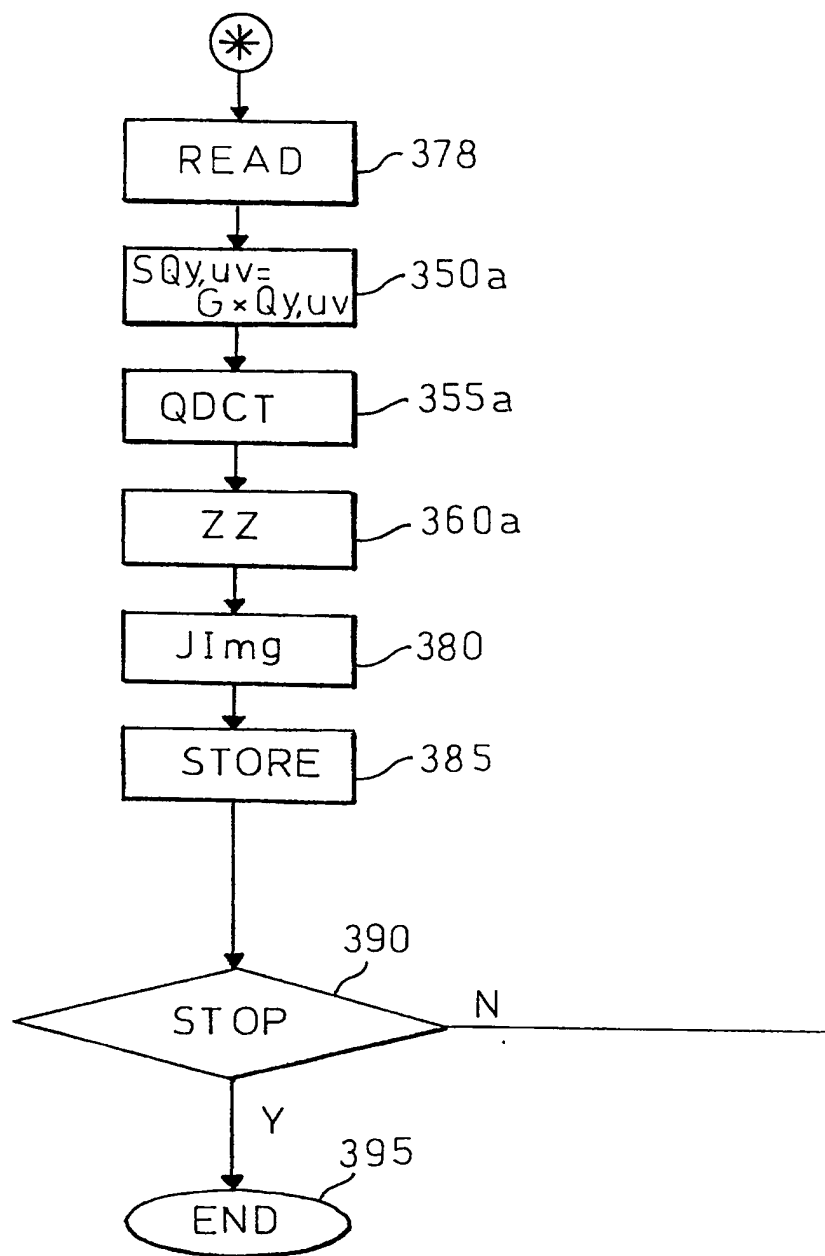


FIG. 3b